

# Pattern recognition systems – Lab 10

## Linear Classifiers and the Perceptron Algorithm

### 1. Objectives

This lab session presents the perceptron learning algorithm for the linear classifier. We will apply gradient descent and stochastic gradient descent procedure to obtain the weight vector for a two-class classification problem.

### 2. Theoretical Background

The goal of classification is to group items that have similar features values into a single class or group. A linear classifier achieves this goal via a discriminant function that is the linear combination of the features.

#### Definitions

Define a training set as the tuple  $(X, Y)$ , where  $X \in M_{n \times m}(R)$  and  $Y$  is a vector  $Y \in M_{n \times 1}(D)$ , where  $D$  is the set of class labels.  $X$  represents the concatenation the feature vectors for each sample from the training set, where each row is an  $m$  dimensional vector representing a sample.  $Y$  is the vector the desired outputs for the classifier. A classifier is a map from the feature space to the class labels:  $f: R^m \rightarrow D$ .

Thus a classifier partitions the feature space into  $|D|$  **decision regions**. The surface separating the classes is called **decision boundary**. If we have only two dimensional feature vectors the decision boundaries are lines or curves. In the following we will discuss binary classifiers. In this case the set of class labels contains exactly two elements. We will denote the labels for classes as  $D = \{-1, 1\}$ .

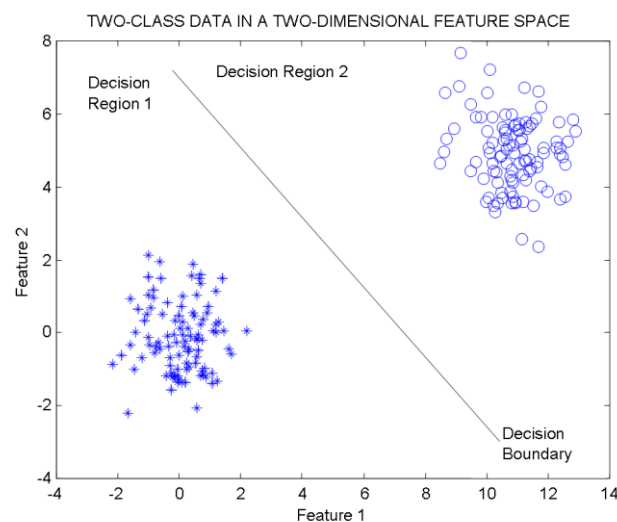


Figure 1. Example of linear classifier on a two-class classification problem. Each sample is characterized by two features.

## 2.1. General form of a linear classifier

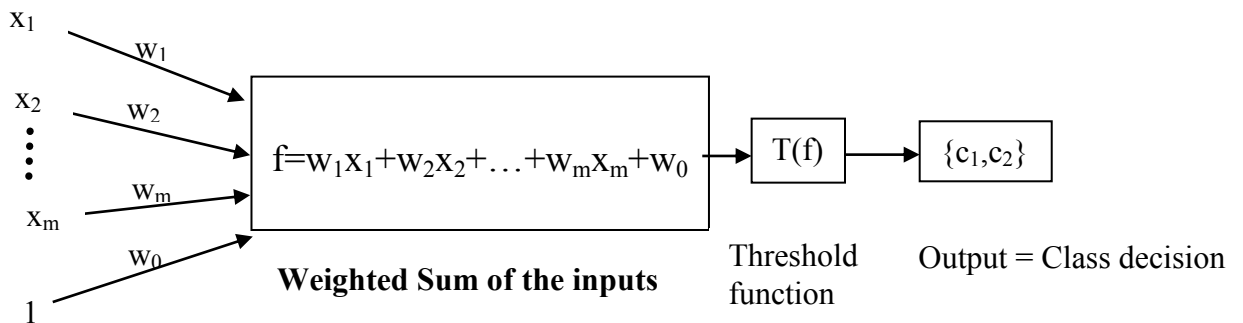
The simplest classifier is a linear classifier. A linear classifier outputs the class labels based on a linear combination of the input features. Considering  $x \in M_{n \times 1}(R)$  as a feature vector we can write the linear decision function as:

$$g(x) = wx^T + w_0 = \sum_{i=1}^n w_i x_i + w_0$$

Where

- $w$  is the **weight vector**
- $w_0$  is the **bias** or the **threshold weight**

A schematic view of the linear classifier is given in the next figure.



For convenience, we will absorb the intercept  $w_0$  by augmenting the feature vector  $x$  with an additional constant dimension (let the bar over a variable denote the augmented version of the vector):

$$wx^T + w_0 = \begin{bmatrix} w_0 & w \end{bmatrix} \begin{bmatrix} 1 \\ x^T \end{bmatrix} = \bar{w}\bar{x}^T$$

A two-category linear classifier (or binary classifier) implements the following decision rule:

$$\begin{cases} \text{if } g(x) > 0 \text{ decide that sample } x \text{ belongs to class } +1 \\ \text{if } g(x) < 0 \text{ decide that sample } x \text{ belongs to class } -1 \end{cases}$$

or

$$\begin{cases} \text{if } w^T x > -w_0 \text{ decide that sample } x \text{ belongs to class } +1 \\ \text{if } w^T x < -w_0 \text{ decide that sample } x \text{ belongs to class } -1 \end{cases}$$

If  $g(x) = 0$ ,  $x$  can ordinarily be assigned to either class.

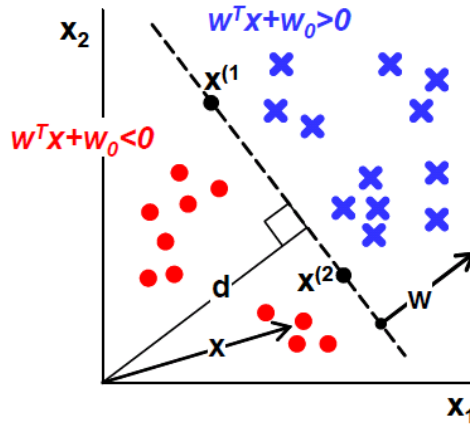


Figure 2. Image for 2D case depicting: linear decision regions (red and blue), decision boundary (dashed line), weight vector ( $w$ ) and bias ( $w_0=d$ ).

## 2.2. Learning algorithms for linear classifiers

We will present two main learning algorithms for linear classifiers. In order to perform learning we transform the task into an optimization problem. For this we define a loss function  $L$ . The loss function applies a penalty for every instance that is classified into the wrong class. The perceptron algorithm adopts the following form for the loss function:

$$L(\bar{w}) = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i \bar{w} \cdot \bar{x}_i^T) = \frac{1}{n} \sum_{i=1}^n L_i(\bar{w})$$

If an instance is classified correctly, no penalty is applied because the second term is negative. In the case of a misclassification the second (positive) term will be added to the function value. The objective now is to find the weights that minimize the loss function.

Gradient descent can be employed to find the global minimum of the loss function. This relies on the idea that a differentiable multivariate function decreases fastest in the opposite direction of the gradient. The update rule according to this observation is:

$$\bar{w}_{k+1} \leftarrow \bar{w}_k - \eta \nabla L(\bar{w}_k)$$

where  $\bar{w}_k$  is the weight vector at time  $k$ ,  $\eta$  is a parameter that controls the step size and is called the learning rate, and  $\nabla L(\bar{w})$  is the gradient vector of the loss function at point  $\bar{w}_k$ . The gradient of the loss function is:

$$\nabla L(\bar{w}) = \frac{1}{n} \sum_{i=1}^n \nabla L_i(\bar{w})$$

$$\nabla L_i(\bar{w}) = \begin{cases} 0, & \text{if } y_i \bar{w} \cdot \bar{x}_i^T > 0 \\ -y_i \bar{x}_i, & \text{otherwise} \end{cases}$$

In the standard gradient descent approach we update the weights after visiting all the training examples. This is also called the batch-update learning algorithm. We can use stochastic gradient descent instead. This entails updating the weights after visiting each training example resulting in the online perceptron learning algorithm. In this case the update rule becomes:

$$\bar{w}_{k+1} \leftarrow \bar{w}_k - \eta \nabla L_i(\bar{w})$$

**Algorithm:**  
**Batch Perceptron**

```
initialize  $\mathbf{w}$ ,  $\eta$ ,  $\theta$ ,  $k=0$ 
for  $k=1:K$ 
   $L = 0$ 
  for  $i=1:n$ 
     $z_i = \sum_{j=0}^d w_j X_{ij}$ 
    if  $z_i \cdot y_i < 0$ 
       $\frac{\partial L}{\partial \mathbf{w}} \leftarrow \frac{\partial L}{\partial \mathbf{w}} - y_i \mathbf{X}_i$ 
       $L \leftarrow L - z_i \cdot y_i$ 
    endif
  endfor
   $L \leftarrow L/n$ 
   $\frac{\partial L}{\partial \mathbf{w}} \leftarrow \frac{\partial L}{\partial \mathbf{w}} / n$ 
  if  $\left\| \frac{\partial L}{\partial \mathbf{w}} \right\| < L_{limit}$ 
    break
   $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial L}{\partial \mathbf{w}}$ 
endfor
```

---

**Algorithm:**  
**Online Perceptron**

```
initialize  $\mathbf{w}$ ,  $\eta$ ,  $\theta$ ,  $k=0$ 
for  $k=1:K$ 
   $E = 0$ 
  for  $i=1:n$ 
     $z_i = \sum_{j=0}^d w_j X_{ij}$ 
    if  $z_i \cdot y_i < 0$ 
       $\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{X}_i y_i$ 
       $E \leftarrow E + 1$ 
    endif
  endfor
   $E \leftarrow E/n$ 
  if  $E < E_{limit}$ 
    break
  endfor
```

---

### 2.3. Two-class two feature linear classifier

In this laboratory session we will find a linear classifier that discriminates between two sets of points. The points in class 1 are colored in red and the points in class 2 are colored in blue.

Each point is described by the color (that denotes the class label) and the two coordinates,  $x_1$  and  $x_2$ .

The augmented weight vector will have the form  $\bar{\mathbf{w}} = [w_0 \ w_1 \ w_2]$ .

The augmented feature vector will be  $\bar{\mathbf{x}} = [1 \ x_1 \ x_2]$ .

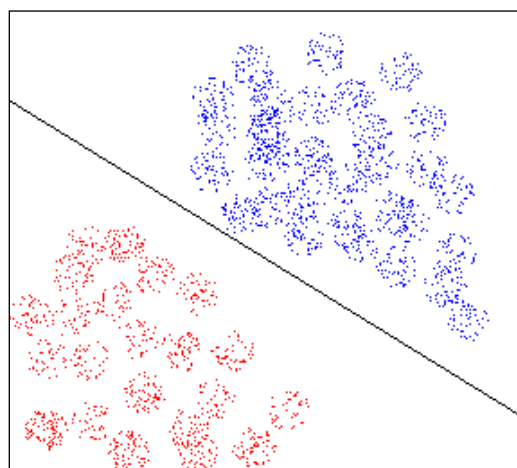


Figure 3. The decision boundary obtained from the perceptron algorithm

### 3. Practical work

1. Read the points from the file *test0\*.bmp* and construct the training set  $(X, Y)$ . Assign the class label +1 to blue points and -1 to red points.
2. Implement and apply the online perceptron algorithm to find the linear classifier that divides the points into two groups. Suggestion for parameters:  $\eta=10^{-4}$ ,  $w_0 = [0.1, 0.1, 0.1]$  and  $E_{limit}=10^{-4}$ .
3. Draw the final decision boundary based on the weight vector  $\mathbf{w}$ .
4. Implement the batch perceptron algorithm and find suitable parameters values. Show the loss function at each step. It must decrease slowly.
5. Visualize the decision boundary at intermediate steps, while the learning algorithm is running.
6. Change the starting values for the weight vector  $\mathbf{w}$ , the learning rate and terminating conditions to observe what happens in each case. Why is setting  $\mathbf{w}$  to the zero vector not a good starting value?

### 4. References

- [1] Richard O. Duda, Peter E. Hart, David G. Stork: *Pattern Classification 2<sup>nd</sup> ed.*
- [2] <http://web.engr.oregonstate.edu/~xfern/classes/cs534/notes/perceptron-4-11.pdf>
- [3] [http://en.wikipedia.org/wiki/Gradient\\_descent](http://en.wikipedia.org/wiki/Gradient_descent)